

Report Documentation Page				Form Approved OMB No. 0704-0188	
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE <b>NOV 2014</b>		2. REPORT TYPE		3. DATES COVERED <b>00-00-2014 to 00-00-2014</b>	
4. TITLE AND SUBTITLE <b>SNUMedinfo at TREC Web Track 2014</b>				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>Seoul National University,Medical Informatics Laboratory,Seoul, Republic of Korea,</b>				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release; distribution unlimited</b>					
13. SUPPLEMENTARY NOTES <b>presented in the proceedings of the Twenty-Third Text REtrieval Conference (TREC 2014) held in Gaithersburg, Maryland, November 19-21, 2014. The conference was co-sponsored by the National Institute of Standards and Technology (NIST) and the Defense Advanced Research Projects Agency (DARPA).</b>					
14. ABSTRACT <b>This paper describes the participation of the SNUMedinfo team at the TREC Web track 2014. This is the first time we participate in the Web track. Rather than applying more sophisticated retrieval method such as learning to rank models, this year we used only baseline retrieval models with spam filtering and pagerank prior.</b>					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT <b>Same as Report (SAR)</b>	18. NUMBER OF PAGES <b>2</b>	19a. NAME OF RESPONSIBLE PERSON
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE <b>unclassified</b>			

# SNUMedinfo at TREC Web track 2014

Sungbin Choi, Jinwook Choi

Medical Informatics Laboratory, Seoul National University, Seoul, Republic of Korea

wakeup06@empas.com, jinchoi@snu.ac.kr

**Abstract.** This paper describes the participation of the SNUMedinfo team at the TREC Web track 2014. This is the first time we participate in the Web track. Rather than applying more sophisticated retrieval method such as learning to rank models, this year we used only baseline retrieval models with spam filtering and pagerank prior.

**Keywords:** Web search, Information retrieval, Sequential dependence model, Spam filtering

## 1. Introduction

In this paper, we describe the methods in participation of the SNUMedinfo team at the TREC Web track 2014. For a detailed task introduction, please see the overview paper of this track.

## 2. Methods

We used sequential dependence model (SDM) [1] as a baseline retrieval model. For the experiment, we used batch query service offered by lemur project website [2]. Clue-Web12-Full dataset is our test corpus. Waterloo spam filter [3] is used to filter out spam documents. Details of our submitted runs can be summarized as following table.

**Table 1. Submitted runs**

RunID	Method description
SNUMedinfo11	SDM
SNUMedinfo12	SDM + Spam filtering (threshold: 50)
SNUMedinfo13	SDM + Spam filtering (threshold: 50) + Pagerank Prior score

*SDM : Sequential dependence model*

Regarding SNUMedinfo13, we used Pagerank Prior [4] scores offered by lemur project website.

### 3. Results

**Table 2. Evaluation results**

RunID	ndcg@20	err@20
SNUMedinfo11	0.2436	0.1386
SNUMedinfo12	0.2698	0.1759
SNUMedinfo13	0.1927	0.1230

### 4. Conclusion

This year, we submitted baseline retrieval model with spam filtering and pagerank prior score. We plan to experiment with more advanced retrieval methods in the next year's participation.

### 5. Acknowledgements

This study was supported by a grant of the Korean Health Technology R&D Project, Ministry of Health & Welfare, Republic of Korea. (No. HI11C1947)

### 6. References

1. Metzler, D. and W.B. Croft, *A Markov random field model for term dependencies*, in *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*. 2005, ACM: Salvador, Brazil. p. 472-479.
2. *The Lemur Project*. [cited 2014 Oct 28]; Available from: <http://www.lemurproject.org>.
3. Cormack, G., M. Smucker, and C.A. Clarke, *Efficient and effective spam filtering and re-ranking for large web datasets*. *Information Retrieval*, 2011. **14**(5): p. 441-465.
4. Page, L., et al., *The PageRank Citation Ranking: Bringing Order to the Web*. 1999, Stanford InfoLab.